# LA-UR-16-26126

Title: Comparison of High Performance Networks: EDR InfiniBand vs. 100Gb RDMA Capable Ethernet

Author(s): Van Wig, Faith Virginia
Kachelmeier, Luke Anthony
Erickson, Kari Natania

Intended for: HPC Mini Showcase

Issued: 2016-08-08

# OUTLINE

- Background
- Overview
- Experimental Setup
- Results
- Difficulties Faced
- Conclusions
- Future Work

# Background Information

- Ethernet

  TCP / IP Stack

- InfiniBand (IB)

  RDMA

  EDR - 4 lanes of 25Gb/s

- RoCE - RDMA over Converged Ethernet



The Seven Layers of OSI

Ethernet → Layer 3 (Network Layer)
InfiniBand → Layer 2 (Data Link Layer)

# OVERVIEW



**TOP500 Interconnect Trends**

http://blog.infinibandta.org/tag/top500/

# PROJECT DESCRIPTION

Compare EDR InfiniBand, RoCE, and native Ethernet:

- Bandwidth Performance

- Latency Performance

- Ease of Configuration (or lack thereof)

OSU Benchmarks, Intel Micro Benchmarks

# Experimental Setup

https://www.attotech.com/products/cables-and-accessories/none/cables-and-accessories/CBL-0310-020

# HARDWARE

- Mellanox ConnectX-4 EDR adapter cards

- Mellanox SB7700 EDR InfiniBand Switch

- Juniper QFX5200 100Gb Ethernet Switch

6

# SERVICES

- CentOS 6 across the nodes

- OpenMPI 1.10 using MXM (Mellanox) libraries

- MLNX OFED drivers and related modules, *v. 3.3-1*

# RESULTS

https://www.healthlottery.co.uk/results/

# Point-to-Point Connected Tests

OSU Benchmarks

Latency (No switch)

LESS IS BETTER

13.55μs

InfiniBand
RoCE
Ethernet

1.814μs

1.752μs

Latency (microseconds)

Packet Size (bytes)

*NOT Log Scaling

10

# Latency (No switch)



LESS IS BETTER

Latency (microseconds) vs Packet Size (bytes)

13.55μs

- InfiniBand
- RoCE
- Ethernet

1.814μs

1.752μs

*NOT Log Scaling

10

Latency (Through switch)

Latency (microseconds)

2.544μs

Δ0.698μs

1.846μs

InfiniBand

RoCE

Packet Size (bytes)

*NOT Log Scaling

12

Bandwidth (No switch)

Bandwidth (No switch)

Bandwidth (No switch)

14

*NOT Log Scaling

**Bandwidth (Through switch)**

Y-axis: Bandwidth (Gigabits) — 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100

X-axis: Packet Size (bytes) — 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096, 8192, 16384, 32768, 65536, 131072, 262144, 5244288, 1048576, 2097152, 4194304

98.7Gb/s

Δ59.6Gb/s

39.1Gb/s

InfiniBand

RoCE

*NOT Log Scaling

15

Message Rate (8 ranks) (No switch)

# Message Rate (8 ranks) (Through switch)

16,096,099 Messages/s

16,010,006 Messages/s

1,064,575 Messages/s

*Uses Log Scaling

- InfiniBand
- RoCE
- Ethernet

**Messages/seconds** (y-axis): 0, 2000000, 4000000, 6000000, 8000000, 10000000, 12000000, 14000000, 16000000, 18000000

**Packet Size (bytes)** (x-axis): 1, 4, 16, 64, 256, 1024, 4096, 16384, 65536, 262144, 1048576, 4194304

*NOT Log Scaling

17

# COLLECTIVE TESTS THROUGH SWITCHES

Intel Micro Benchmarks

All-to-all

Latency (microseconds) vs Packet Size (bytes)

29.19s
1.495s
0.333s

InfiniBand
RoCE
Ethernet

*NOT Log Scaling

*Uses Log Scaling

19

All-to-all-v

All-reduce

*NOT Log Scaling

*Uses Log Scaling

21

# Difficulties Faced

http://fidelisspies.blogspot.com/2014/03/day-212-can-i-live-without-frustration.html

# InfiniBand Setup Issues

- Older version of InfiniBand hardware better initial performance

- Collective tests hung indefinitely



**InfiniBand Latency FDR vs EDR**

NEWER HARDWARE

OLDER HARDWARE

FDR

EDR

# Impact of InfiniBand Optimization

One-time driver update:

2.2-1 → 3.3-1

Significant latency improvement!

Drivers will be updated throughout HPC division



Latency Before and After Driver Update

- EDR Before
- EDR After

Latency (microseconds) vs Packet Size (bytes)

# Ethernet Setup Issues

- Initial iperf tests (with 6 sockets) were only 60Gb/s
  - Eventually reached 87Gb/s through *lots* of configuration

```
# From Mellanox Tuning Guide and subsequent testing
net.ipv4.tcp_timestamps = 0
net.ipv4.tcp_sack = 1
net.core.netdev_max_backlog = 250000
net.core.rmem_max = 134217728
net.core.wmem_max = 134217728
net.core.rmem_default = 4194304
net.core.wmem_default = 4194304
net.core.optmem_max = 4194304
net.ipv4.tcp_rmem = 4096 87380 134217728
net.ipv4.tcp_wmem = 4096 87380 134217728
net.ipv4.tcp_low_latency = 1
net.ipv4.tcp_adv_win_scale = 1
net.ipv4.tcp_mtu_probing = 1
```

```
ethtool -K eth4 lro on
ethtool --set-priv-flags eth4 hw_lro on
```

- OpenMPI gives <20Gb/s for native Ethernet
  - Potentially lack of socket optimization

# Juniper Setup Issues

- Juniper switch reported Mellanox passive 100Gb copper cables as 40Gb/s

- EDR cards didn't acknowledge Juniper passive optics whatsoever

- Juniper switch didn't acknowledge Mellanox active optics

# MELLANOX CARD COMPATIBILITIES

**Table 12 - Tested 10Gb/s/EDR Switches**

| Speed | Switch Silicon | OPN # / Name | Description | Vendor |
|---|---|---|---|---|
| EDR | Switch-IB | SB7790-EB2F | 36-port EDR 100Gb/s InfiniBand Switch Systems | Mellanox |
| EDR | Switch-IB 2 | SB7800-ES2R | 36-port Non-blocking Managed EDR 100Gb/s InfiniBand Smart Switch | Mellanox |
| 100GbE | Spectrum | SN2700-CS2R | 32-port Non-blocking 100GbE Open Ethernet Spine Switch System | Mellanox |
| 100Gb/s | N/A | C3232C | High-Density, 100 Gigabit Ethernet Switch | Cisco |

# CONCLUSION

http://blog.friendfriend.net/2014/05/10/54-incompatible-plug/

# Conclusion

- Good potential, as shown with direct results.
- 100Gb technology is too new, not standardized
  - Deployment Effort is complex for both options
- Different companies not necessarily compatible
- If you want 100Gb/s, get it all from one place.

# GREATER IMPACT

- Mellanox OFED driver updates throughout HPC Division
- Juniper and Mellanox will begin to collaborate with hardware compatibilities
  - Because of us, they have a place to start

# FUTURE WORK

- Comparison with Intel's OmniPath
- Collective tests after technologies are compatible
- Work with OpenMPI community to optimize Ethernet performance
- NFS over RDMA
- TCP/IP over InfiniBand

# ACKNOWLEDGEMENTS

Thank You!
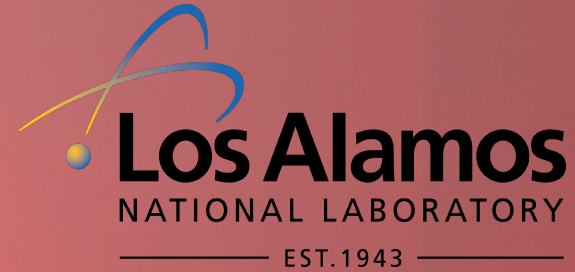
Los Alamos
NATIONAL LABORATORY
— EST.1943 —

Our mentors:

Susan Coulter

Howard Pritchard

Mellanox
TECHNOLOGIES

Aleksey Senin

Martijn Van Breugel

JUNIPER
NETWORKS

Erik DeHaas

Matthew Herzog

Our instructors:

Jarrett Crews

Eric Holm

32

# QUESTIONS?

Faith Van Wig                                              fvvnr5@mst.edu

Luke Kachelmeier                                   lkachelmeier@unm.edu

Kari Erickson                        kari.erickson@student.nmt.edu

# Backup Slides

```
# Kernel sysctl configuration file for Red Hat Linux
#
# For binary values, 0 is disabled, 1 is enabled.  See sysctl(8) and
# sysctl.conf(5) for more details.
#
# Use '/sbin/sysctl -a' to list all possible parameters.

# Controls IP packet forwarding
net.ipv4.ip_forward = 0

# Controls source route verification
net.ipv4.conf.default.rp_filter = 1

# Do not accept source routing
net.ipv4.conf.default.accept_source_route = 0

# Controls the System Request debugging functionality of the kernel
kernel.sysrq = 0

# Controls whether core dumps will append the PID to the core filename.
# Useful for debugging multi-threaded applications.
kernel.core_uses_pid = 1

# Controls the use of TCP syncookies
net.ipv4.tcp_syncookies = 1

# Controls the default maxmimum size of a mesage queue
kernel.msgmnb = 65536

# Controls the maximum size of a message, in bytes
kernel.msgmax = 65536

# Controls the maximum shared segment size, in bytes
kernel.shmmax = 68719476736

# Controls the maximum number of shared memory segments, in pages
kernel.shmall = 4294967296

# From Mellanox Tuning Guide and subsequent testing
net.ipv4.tcp_timestamps = 0
net.ipv4.tcp_sack = 1
net.core.netdev_max_backlog = 250000
net.core.rmem_max = 134217728
net.core.wmem_max = 134217728
net.core.rmem_default = 4194304
net.core.wmem_default = 4194304
net.core.optmem_max = 4194304
net.ipv4.tcp_rmem = 4096 87380 134217728
net.ipv4.tcp_wmem = 4096 87380 134217728
net.ipv4.tcp_low_latency = 1
net.ipv4.tcp_adv_win_scale = 1
net.ipv4.tcp_mtu_probing = 1
```

ifcfg-eth4

```
DEVICE=eth4
BOOTPROTO=dhcp
PEERDNS=no
ONBOOT=yes
GATEWAY=192.168.1.1
MTU=9000
```

rc.local (Onboot)

```
ethtool -K eth4 lro on
ethtool --set-priv-flags eth4 hw_lro on
```

whenever you run OMPI:

```
export OMPI_MCA_tcp_max_send_size=2097152
export OMPI_MCA_btl_tcp_if_include=eth4
export OMPI_MCA_btl_tcp_eager_limit=8388608
export OMPI_MCA_btl_tcp_max_send_size=2097152
export OMPI_MCA_btl_tcp_links=1
export OMPI_MCA_btl_tcp_sendbuf=4194304
export OMPI_MCA_btl_tcp_recvbuf=4194304
```